# Computing at CDF

**Mark Neubauer**
*Massachusetts Institute of Technology*
for the CDF Collaboration

- Introduction
- Computing requirements
- Central Analysis Facility
- Data Handling
- Toward the Grid
- Conclusions

# CDF in a Nutshell

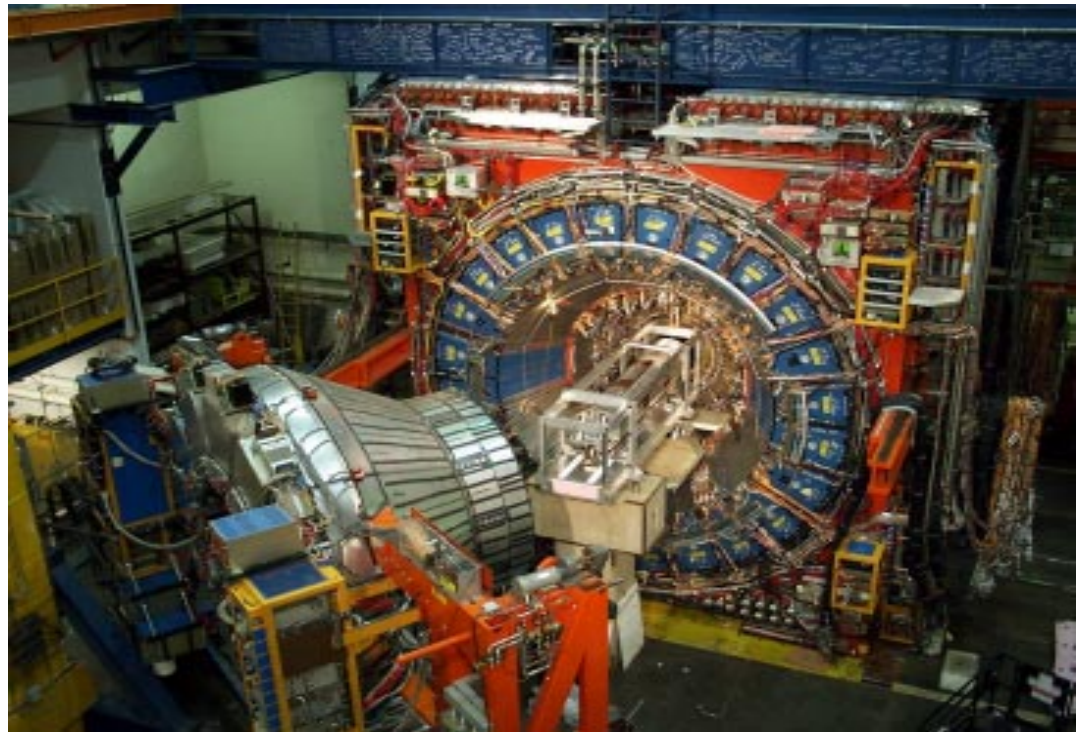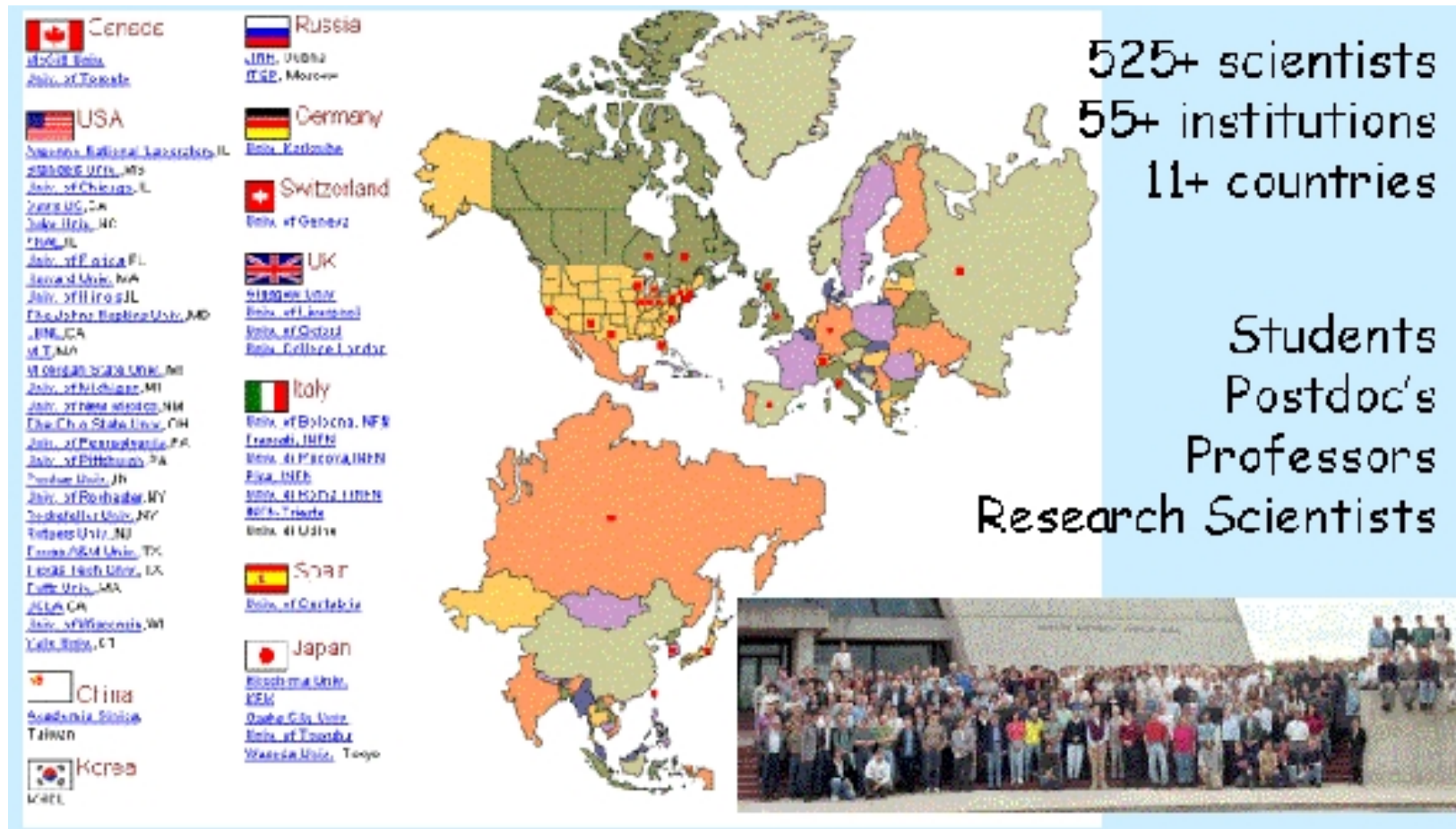- **CDF + D0 experiments analyze p$\bar{\text{p}}$ collisions from Tevatron at Fermilab**
- **Tevatron highest energy collider in world ($\sqrt{s} = 2$ TeV) until LHC**
- **Run I (1992-1996) huge success $\rightarrow$ 200+ papers (t quark discovery, ...)**
- **Run II (March 2001-) upgrades for luminosity ($\times$10) + energy (~10%$\uparrow$)**
   - $\rightarrow$ expect integrated luminosity 20$\times$ (Run IIa) and 150$\times$ (Run IIb) of Run I
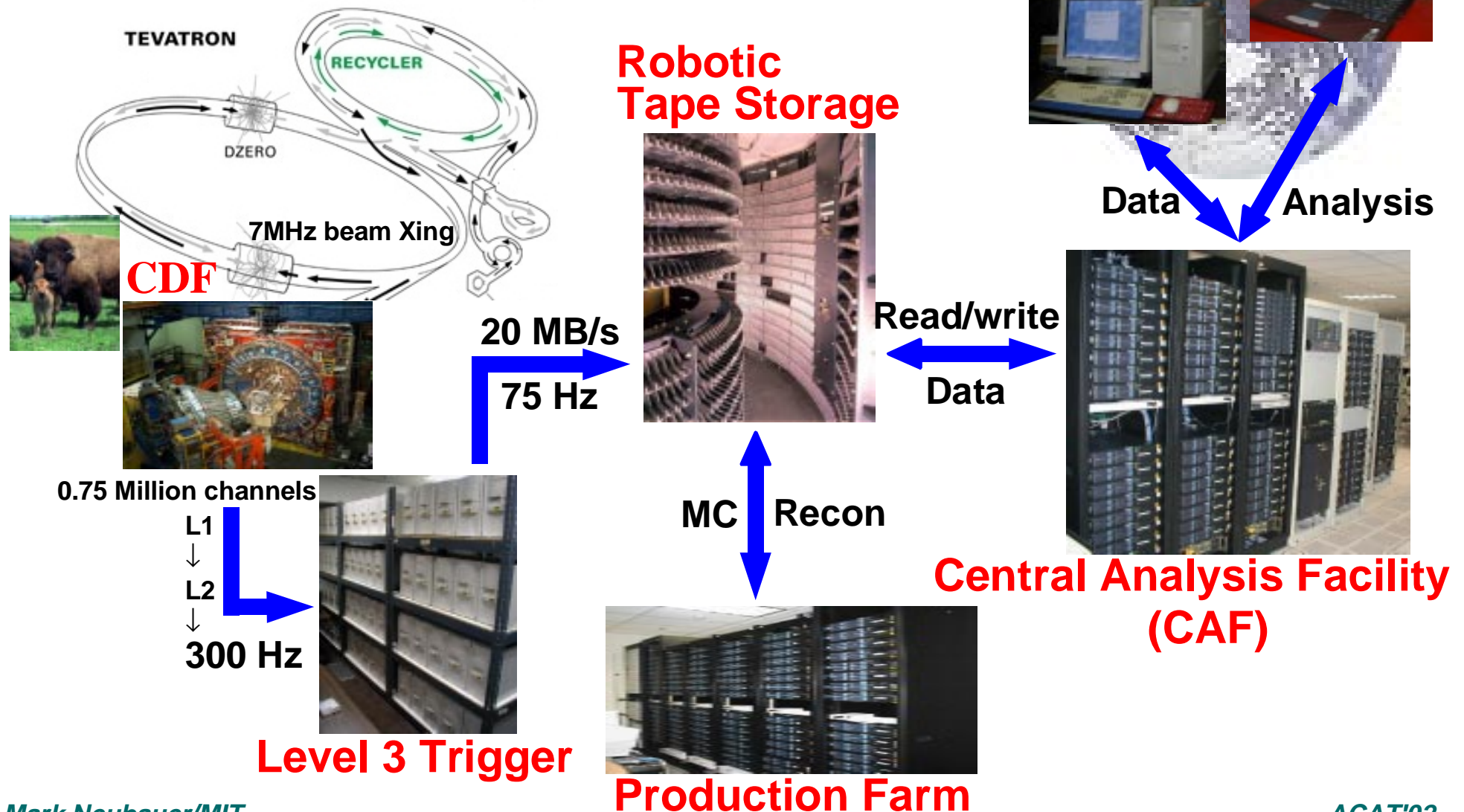
## Run II physics goals:



- **Search for Higgs boson**
- **Top quark properties** (m$_t$, $\sigma_{tot}$, ...)
- **Electroweak** (m$_W$, $\Gamma_W$, ZZ$\gamma$, ...)
- **Search for new physics** (e.g. SUSY)
- **QCD at large Q$^2$** (jets, $\alpha_s$, ...)
- **CKM tests in *b* hadron decays**

# CDF RunII Collaboration



525+ scientists
55+ institutions
11+ countries

Students
Postdoc's
Professors
Research Scientists

**Goal**: **Provide computing resources for 200+ collaborators simultaneously doing analysis per day!**
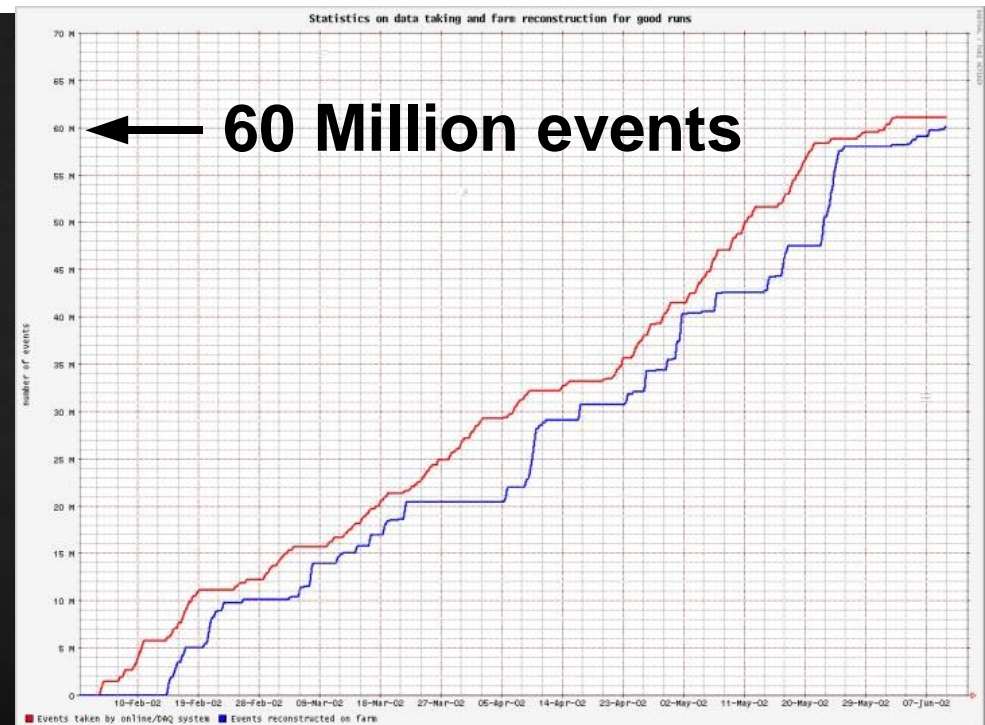
# CDF DAQ/Analysis Flow



TEVATRON

RECYCLER

DZERO

7MHz beam Xing

**CDF**

**User Desktops**

**Robotic Tape Storage**

**Data** **Analysis**

0.75 Million channels

L1
↓
L2
↓
**300 Hz**

**20 MB/s**

**75 Hz**

**Read/write**

**Data**

**MC** **Recon**

**Level 3 Trigger**

**Production Farm**

**Central Analysis Facility (CAF)**

# Reconstruction Farms

**Data reconstruction + validation, Monte Carlo generation**

**154 dual P3's (equivalent to 244 1 Ghz machines)**

**Job management:**

➢ **Batch system → FBSNG developed at FNAL**

➢ **Single executable, validated offline**
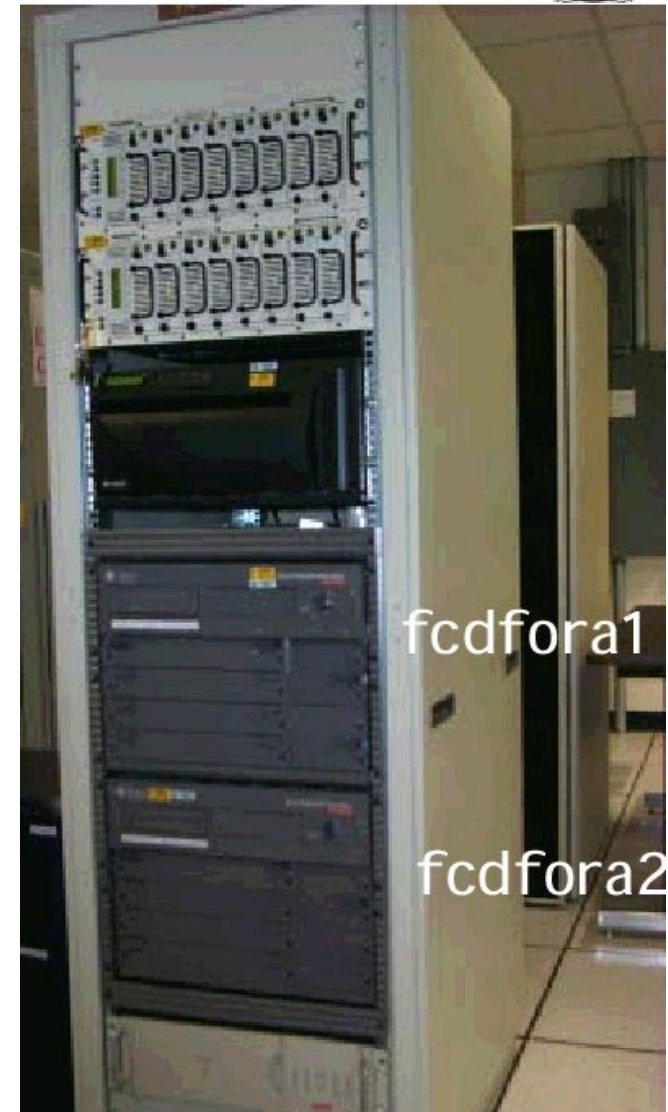


**← 60 Million events**

# Database Usage at CDF

**Oracle DB: Metadata + Calibrations**

## DB Hardware:

- **2 Sun E4500**

## Presently evaluating:

- **Oracle on Linux**
- **MySQL**
- **Replication to remote sites**

fcdfora1

fcdfora2

# **Data/Software Characteristics**

## Data Characteristics:

- ➢ **Root I/O as persistent (raw) data format**
- ➢ **Raw data size: ~250 kB/event**
- ➢ **Reconstructed data (PAD) format: 50-100 kB/event**
- ➢ **Typical ntuple size (stntuple): 5-10 kB/event**
- ➢ **Typical RunIIa secondary dataset size: $10^7$ events**

## Analysis Software:

- ➢ **Typical analysis jobs run @ 10 Hz on 1 GHz P3**
  - $\rightarrow$ **few MB/sec**
- ➢ **CPU rather than I/O bound (FastEthernet)**

# Computing Requirements



**Requirements set by goal:**
200 simultaneous users to analyze secondary data set ($10^7$ evts) in a day

**Need ~700 TB of disk and ~5 THz of CPU by end of FY'05:**
$\rightarrow$ **need lots of disk** $\rightarrow$ **need cheap disk** $\rightarrow$ **IDE Raid**
$\rightarrow$ **need lots of CPU** $\rightarrow$ **commodity CPU** $\rightarrow$ **dual Intel/AMD**

# Past CAF Computing Model



Large SMP (128 processor SGI)
Expensive disks (FiberChannel/SCSI)

Analysis Code Development
Analysis Job Debugging
Interactive Analysis Jobs
Batch Jobs
"Other" Usage

fcdfsgi2

**Very expensive to expand and maintain**

**Bottom line:**
   **Not enough 'bang for the buck'**

# CAF Implementation



**Users are able to:**
- submit jobs
- monitor job progress
- retrieve output

**from 'any desktop' in the world**

# CAF Milestones

➢ **Start of CAF design**    **11/01**

➢ **CAF prototype (protoCAF) assembled**    **2/25/02**

➢ **Fully-functional prototype system (>99% job success)**    **3/6/02**

➢ **ProtoCAF integrated into Stage1 system**    **4/25/02**

➢ **Production Stage1 CAF for collaboration**    **5/30/02**

**Design → Production system in 6 months!**

**ProtoCAF**

**Stage1**

# CAF Stage 1 Hardware



Code Server

File Servers

Worker Nodes

Linux 8-ways
**(interactive)**

# Stage 1 Hardware: Workers

**Workers** (**132 CPUs**, 1U+2U rackmount):

16 2U Dual Athelon 1.6GHz / 512MB RAM
50 1U/2U Dual P3 1.26GHz / 2GB RAM
FE (11 MB/s) / 80GB job scratch each

# Stage 1 Hardware: Servers

**Servers (35TB total, 16 4U rackmount):**

2.2TB useable IDE RAID50 hot-swap
Dual P3 1.4GHz / 2GB RAM
SysKonnect 9843 Gigabt Ethernet card

# File Server Performance



Local disk reads

200 MB/s

60 MB/s

Remote reads from CAF file server

70 MB/s

**Server/Client Performance**: Up to **200MB/s local reads, 70 MB/s NFS**

**Data Integrity tests**: md5sum of local reads under heavy load
BER $< 2\times10^{-14}$ (Maxtor claims $< 1$ error / $10^{14}$ bits read)

**Cooling tests**: Temp profile of disks w/ IR gun after extended disk thrashing

*Mark Neubauer/MIT*                                                               *ACAT'02*

# CAF Software

**Design goal:**

> **Give users access to CAF resources**
> - ➢ **CPU**
> - ➢ **scratch disk**
> - ➢ **data handling system**
>
> **from their desktops anywhere in the world**

**Design constraints/desirables:**

> **Fermilab computing security policy** $\rightarrow$ **kerberos!**
>
> **Job scheduling** $\rightarrow$ **proven batch system, configurable,**
> **fair share capability, local support** $\rightarrow$ **FBSNG (FNAL-CD)**
>
> **Adminstrative ease** $\rightarrow$ **no user accounts**
> $\rightarrow$ **non-interactive batch, jobs run as single 'cafuser' user**
>
> **User identity** $\rightarrow$ **unique privileges for batch jobs, disk space**

# User Access to CAF

## Job Related:

- Submit jobs
- Check progress of job
- Kill a job

## Remote file system access:

- 'ls' in job's 'relative path'
- 'ls' in a CAF node's absolute path
- 'tail' of any file in job's 'relative path'
- Get full file listing based on metadata

# CAF Software

# CAF User Interface

➢ **Compile, build, debug analysis job on 'desktop'**

section integer range



➢ **Fill in appropriate fields & submit job**

output destination

user exe+tcl directory

➢ **Retrieve output using kerberized FTP tools ... or write output directly to 'desktop'!**

# Web Monitoring of User Queues

**Each user a different queue**

**Process type** for job length

- **test**: 5 mins
- **short**: 2 hrs
- **medium**: 6 hrs
- **long**: 2 days

**This example:**

1 job → 11 sections

(+ 1 additional section automatic for job cleanup)



FBSNG on the web
Farm:        CAF
Time:        Thu May 23 02:32:41 2002
Report:      List of queues

Queues | Jobs | Nodes | Process Types

User Monitor

| Name | Status | Default Process Type | Share | Prio | Waiting | Ready | Running | Total |
|------|--------|----------------------|-------|------|---------|-------|---------|-------|
| akorn | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| amitl | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| anikeev | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| belforte | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| msmartin | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| msn | OK | short | 1.00 | 0 | 1 | 0 | 11 | 12 |
| pauly | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| paus | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| ratnikov | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| rescigno | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| semeria | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| sfiligoi | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| sgromoll | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| shepard | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| sidoti | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| spezziga | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| test | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| thkim | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| thom | OK | short | 1.00 | 0 | 1 | 0 | 1 | 2 |

*Mark Neubauer/MIT*

*ACAT'02*

# Monitoring jobs in your queue



*Mark Neubauer/MIT*

*ACAT'02*

# Monitoring sections of your job

# CAF Utilization



**CAF global status**

ptype | Average % | Current %
--- | --- | ---
Short | 69.0 | 0.0
Medium | 3.5 | 1.5
Long | 2.8 | 2.2
All processes | 75.2 | 3.7

Updated: Jun 18 12:20:01 2002

Built using RRDTool

**Active queues ( last 24h )**

Updated: Jun 18 12:20:03 2002

schuster | test | fkw | anikeev | tomohiro | gotra | akorn | castro
casarsa | cjl | jdlee | zyhanv | ikfuric | lys | nigmanov | gpope
jmuelmen | rescigno | dbstarr | ikrav | daronco | thom

Built using RRDTool

## Status summary

| | Short | Medium | Long | All Types |
| --- | --- | --- | --- | --- |
| Running sections | 0 | 2 | 3 | 5 |
| Pending sections | 0 | 0 | 0 | 0 |
| Waiting time [hh:mm] (24h average): | | | | |
| per job | 0:04 | 0:26 | 0:00 | 0:15 |
| per section | 2:12 | 0:52 | 0:00 | 1:32 |
| Running time [hh:mm] (24h average) | 0:20 | 0:35 | 0:00 | 0:27 |

Updated: Jun 18 12:20:02 2002

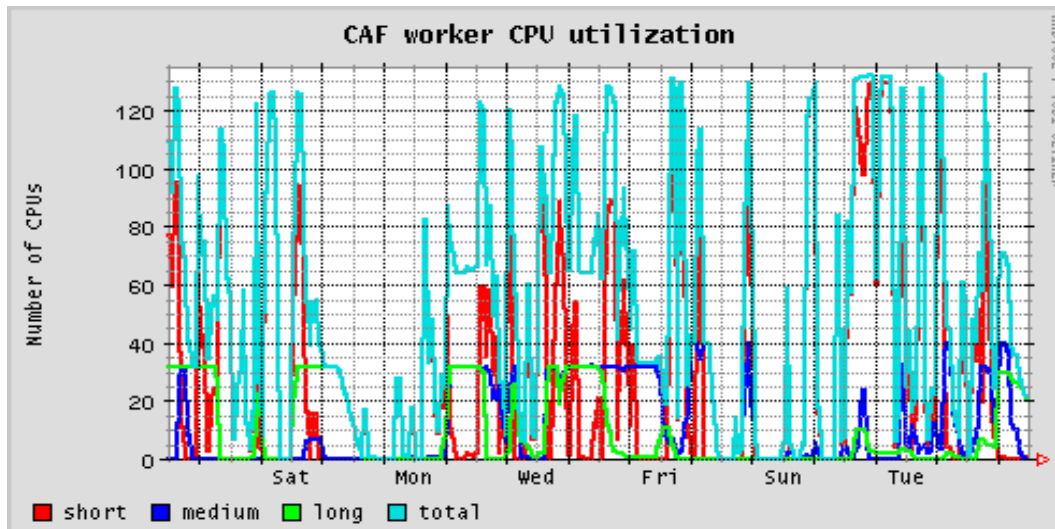## CAF in active use by CDF collaboration

- 120 CAF Users (queues) to date
- 2-5 new users per day
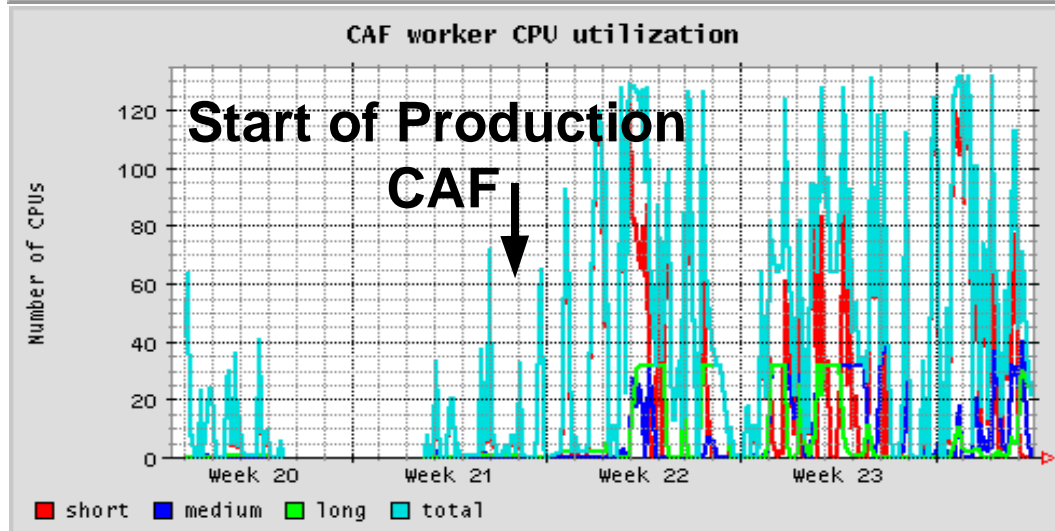- Several dozen simultaneous users in a typical 24 hr period

# CAF Utilization



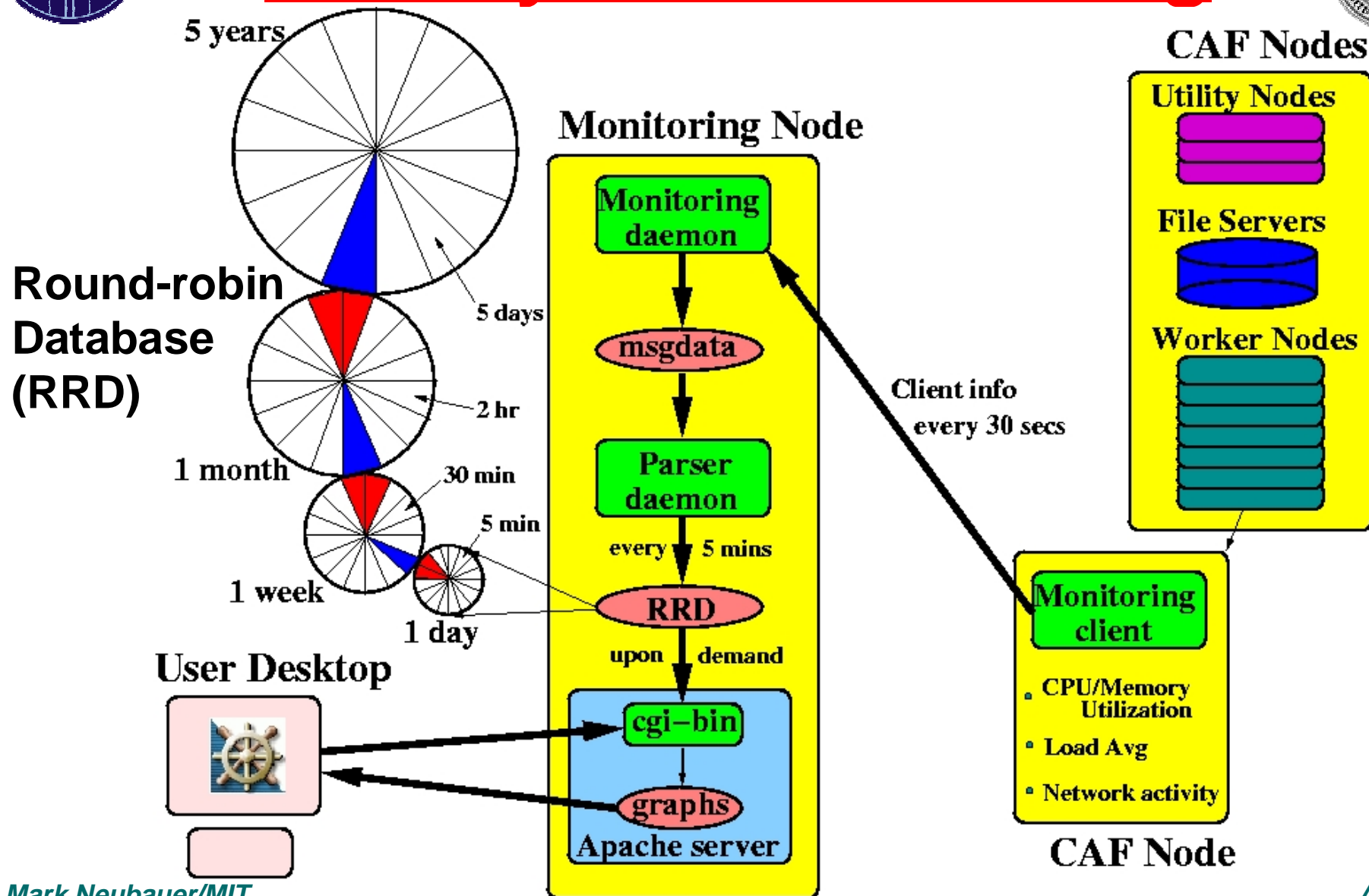CAF utilization steadily rising since opened to collaboration

Providing 10-fold increase in analysis resources for summer physics conferences

Need for more CPU on the horizon
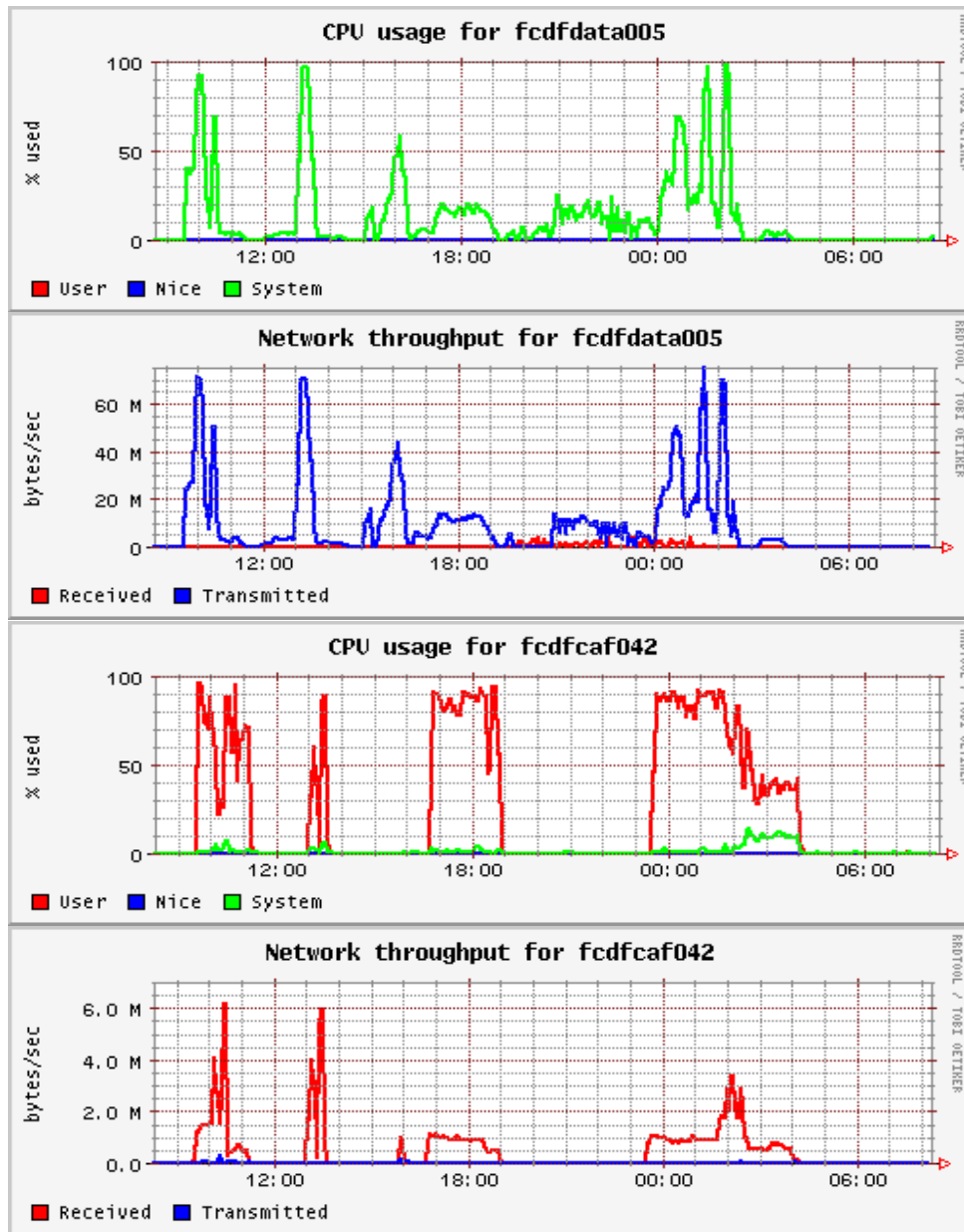
# CAF System Monitoring

# System Monitoring

## 2 TB File Server

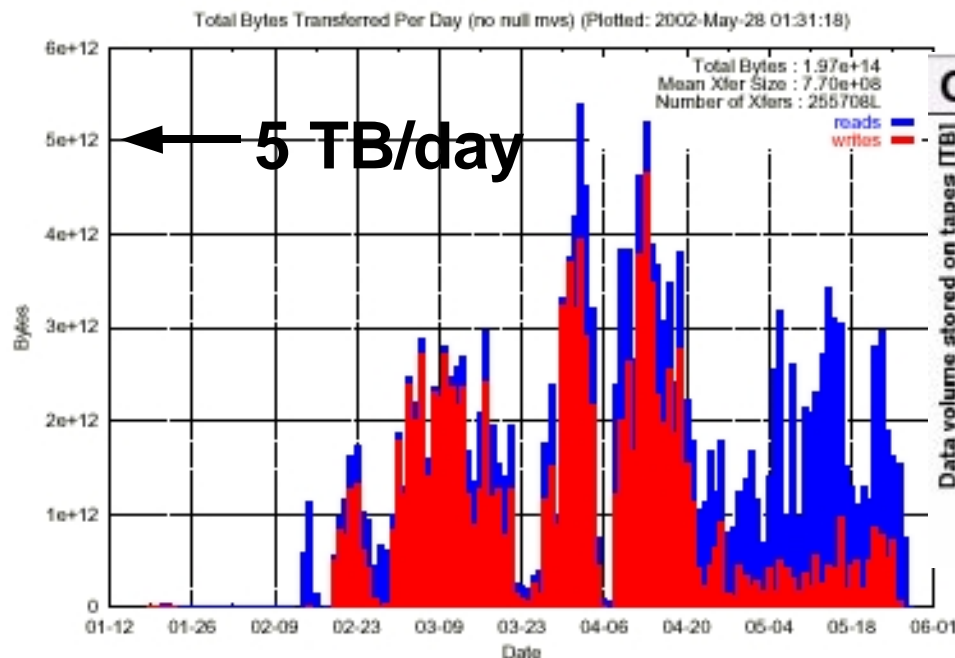**Data transfers CPU limited**

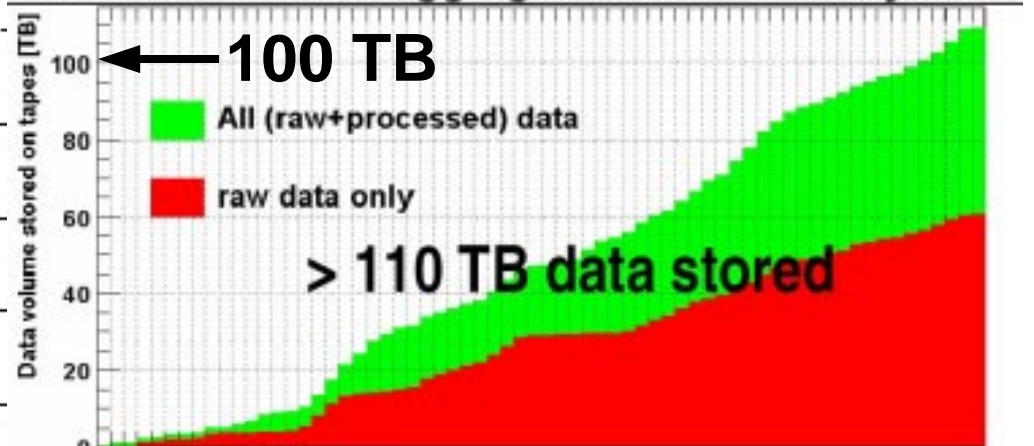**Analysis Jobs CPU bound**

## Worker Node

# Data Handling

**Data archived using STK 9940 drives and tape robot**

**Enstore: Network-attached tape system developed at FNAL**
$\rightarrow$ **provides interface layer for staging data from tape**

# Data Handling

**Dcache $\rightarrow$ network-attached disk cache from DESY**

- Front-end disk cache for Enstore (read and write disk pools)
- Currently in $\beta$ testing $\rightarrow$ working toward production use in CDF
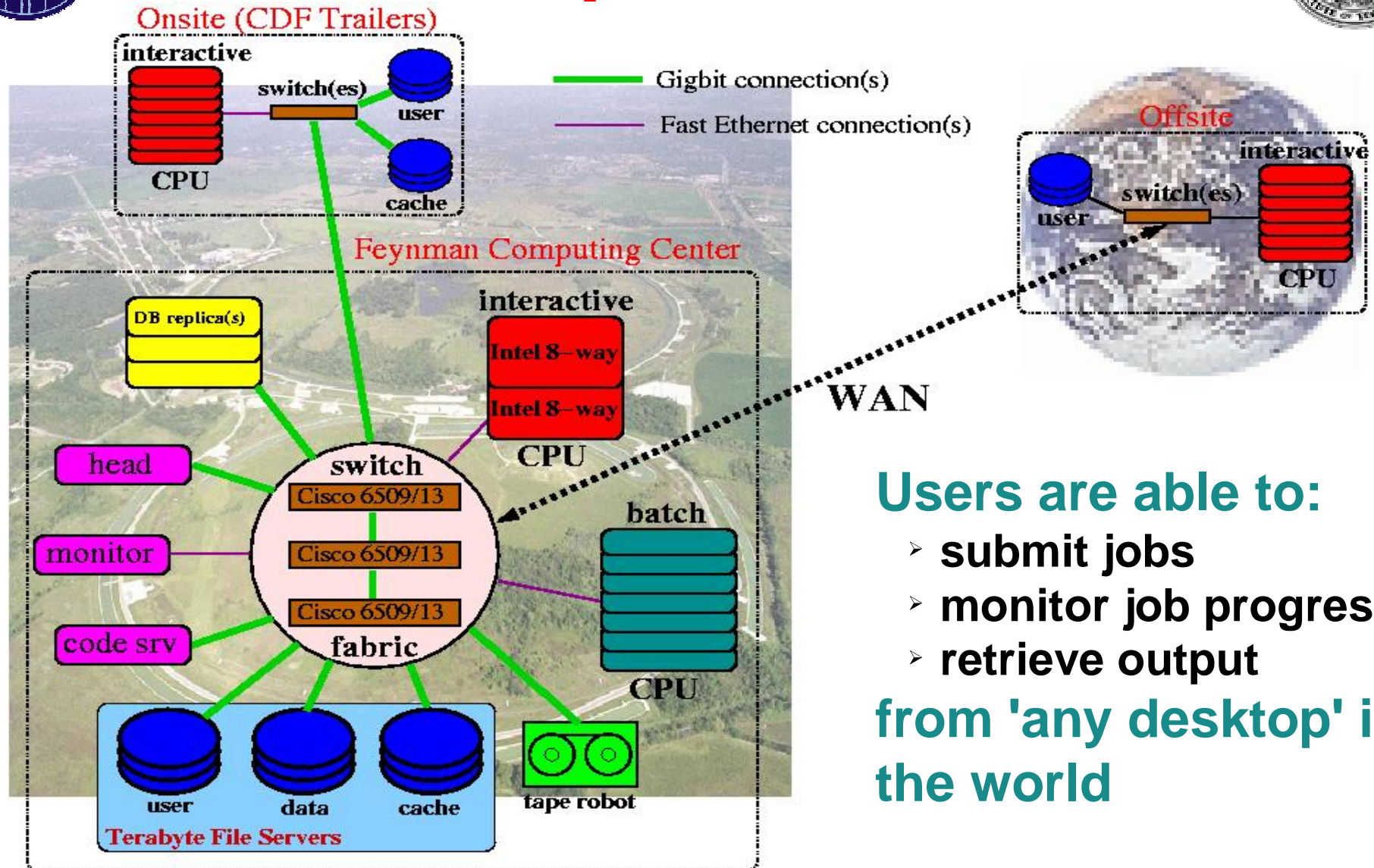
**SAM $\rightarrow$ framework for global data handling/distribution**

- Jointly developed by FNAL Computing division and D0
- Works with Enstore and CDF analysis software framework
- Currently under evaluation for use in CDF data distribution

$\rightarrow$ **see Igor Terekhov's talk**

# CAF Implementation
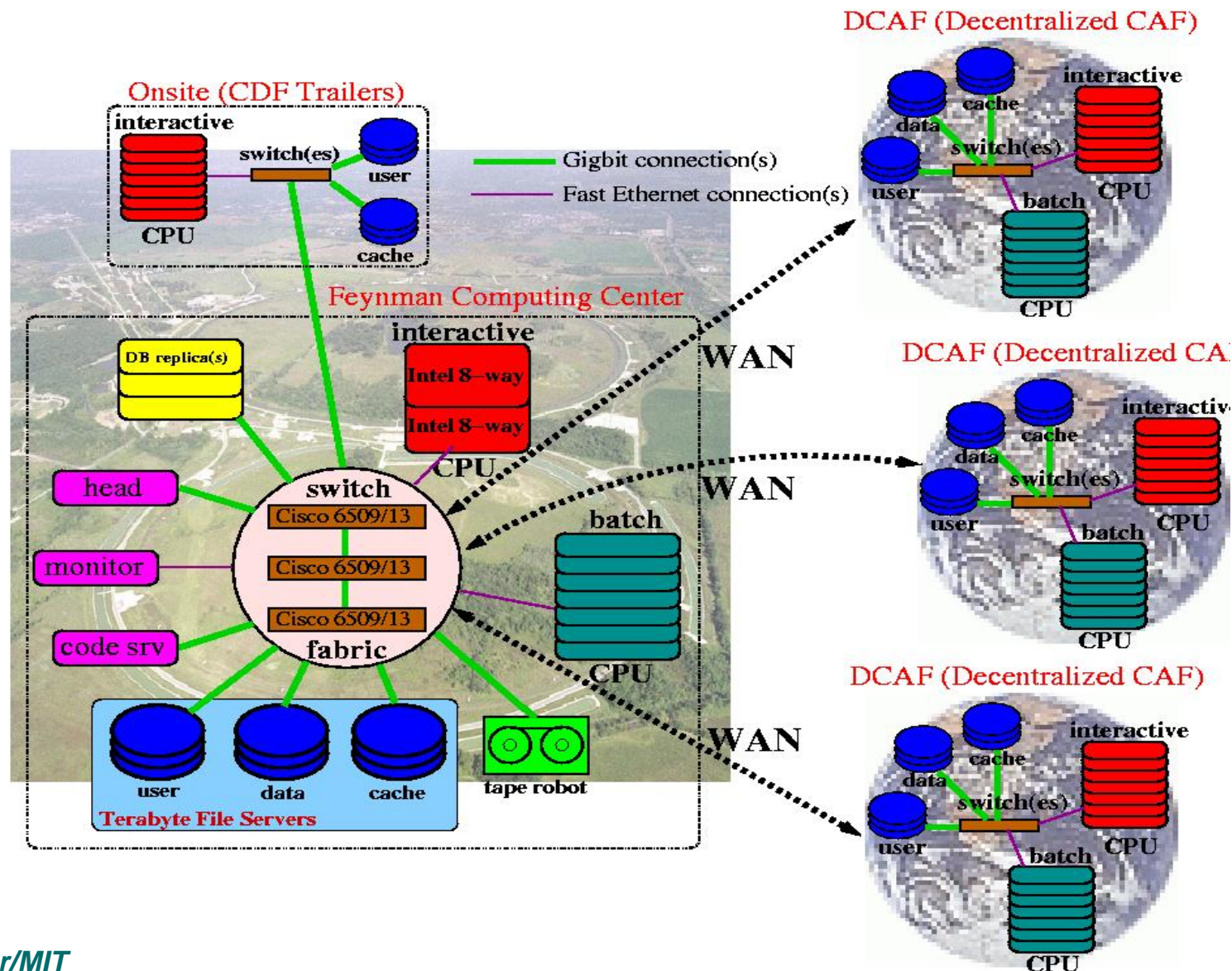


Users are able to:
- submit jobs
- monitor job progress
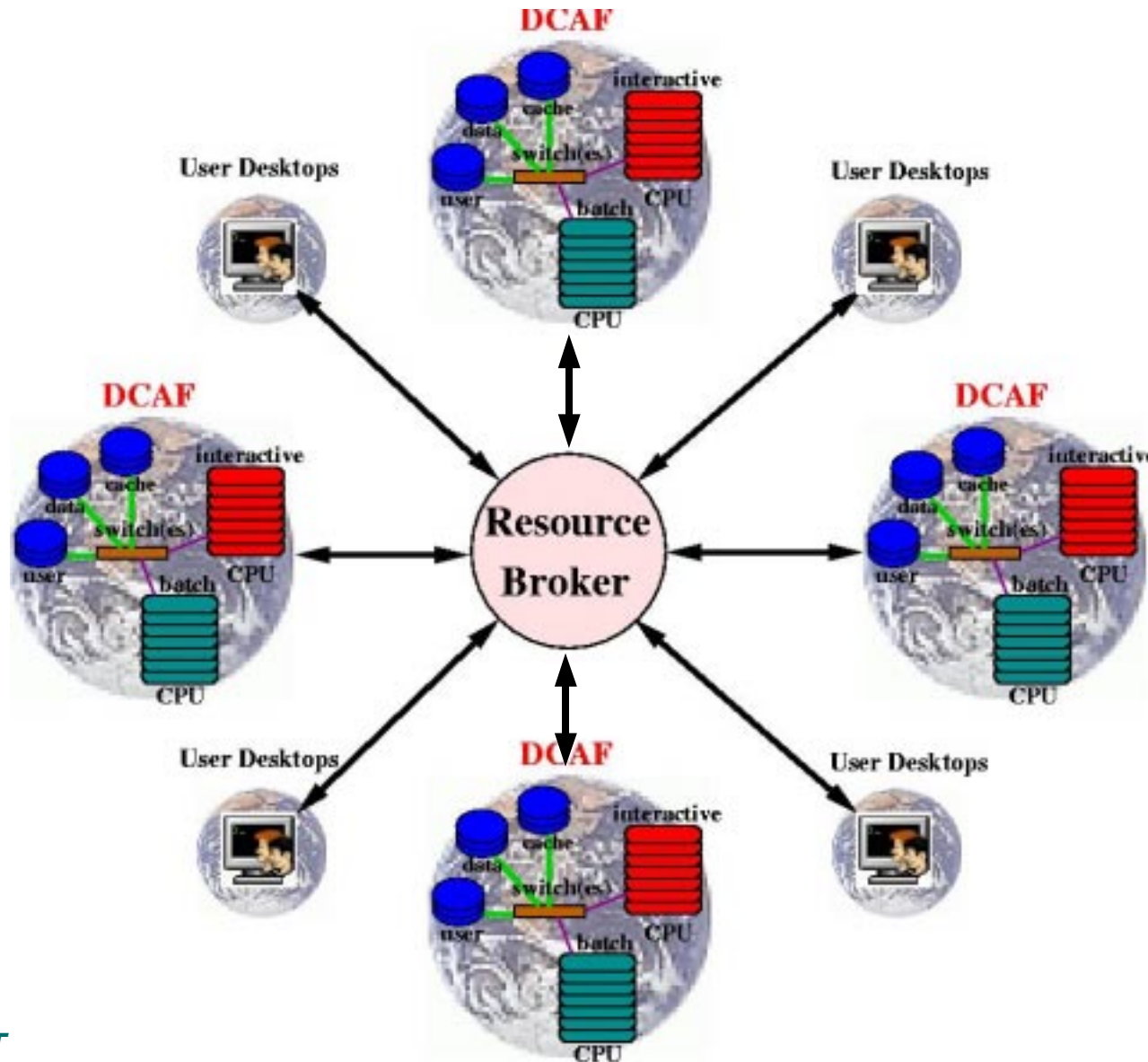- retrieve output

from 'any desktop' in the world

# Toward the Grid

# Peer-to-Farm Paradigm

# Brokering scheme

**Minimize job execution time:**
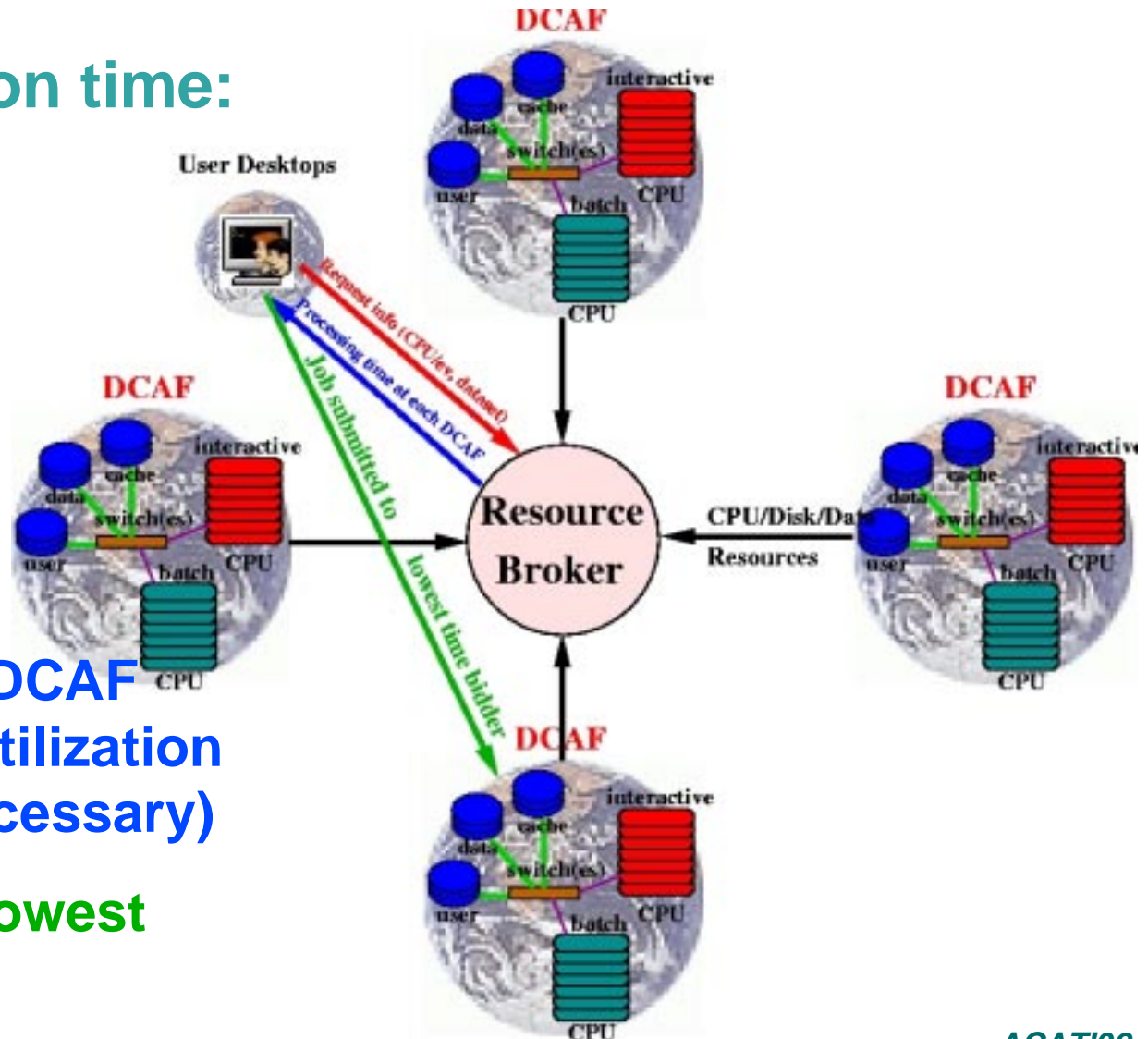
**DCAFs update broker**
- ➢ **CPU/disk utilization**
- ➢ **Local data**

**User generates request**
- ➢ **CPU time/event**
- ➢ **Metadata ID (dataset)**

**Execution time on each DCAF**
- ➢ **CPU+I/O resources+utilization**
- ➢ **Data movement (if necessary)**

**Job goes to DCAF with lowest 'bid'**

# Summary/Conclusions

**Distributed Peer-to-Farm Computing Model**

**Production system under heavy use:**
- ➢ **Single farm at FNAL**
- ➢ **Many peers all over the world**
  - **100+ total users**
  - **100+ simultaneous jobs**
  - **Regularly up to 800 jobs per user queued**

**Future development:**
- ➢ **Extend data handling**
- ➢ **Multi-farm brokering**
- ➢ **Scale system by O(10)**